

1 CAHILL GORDON & REINDEL LLP  
JOEL KURTZBERG (admitted *pro hac vice*)  
2 FLOYD ABRAMS (admitted *pro hac vice*)  
JASON ROZBRUCH (admitted *pro hac vice*)  
3 LISA J. COLE (admitted *pro hac vice*)  
32 Old Slip  
New York, New York 10005  
4 Telephone: 212-701-3120  
Facsimile: 212-269-5420  
5 jkurtzberg@cahill.com

6 DOWNEY BRAND LLP  
WILLIAM R. WARNE (Bar No. 141280)  
bwarne@downeybrand.com  
7 MEGHAN M. BAKER (Bar No. 243765)  
mbaker@downeybrand.com  
8 621 Capitol Mall, 18th Floor  
Sacramento, CA 95814  
9 Telephone: 916-444-1000  
Facsimile: 916-520-5910

10 *Attorneys for Plaintiff X Corp.*

11  
12 UNITED STATES DISTRICT COURT  
13 EASTERN DISTRICT OF CALIFORNIA  
14 SACRAMENTO DIVISION

15  
16 X CORP.,

Plaintiff,

17 v.

18 ROBERT A. BONTA, Attorney  
General of California, in his  
19 official capacity,

20 Defendant.

No. 2:23-cv-01939-WBS-AC

**AFFIDAVIT OF [REDACTED]  
IN SUPPORT OF X CORP.'S MOTION  
FOR PRELIMINARY INJUNCTION**

1 [REDACTED], being duly sworn, deposes and states as follows:

2 1. I am a [REDACTED] on the Trust and Safety Team at X Corp,  
3 covering product, policy, and operations. In this role, I assist  
4 with the creation, implementation, and enforcement of policies –  
5 including content moderation policies – on the X platform.

6 2. I am submitting this affidavit in support of Plaintiff's  
7 Motion for Preliminary Injunction. I have personal knowledge of  
8 the facts set forth herein, unless otherwise noted. If called upon  
9 as a witness, I could and would competently testify to those facts.

10 **I. Applicability of AB 587**

11 3. I have reviewed and am familiar with the law in this  
12 case: California's AB 587. AB 587 applies to X Corp. and its X  
13 platform.

14 4. X Corp. is a "social media company," as defined by  
15 § 22675(d), as it is a "person or entity that owns or operates a  
16 social media platform," as defined by the statute. X Corp. owns  
17 and operates X, which is a social media platform, as defined by AB  
18 587.

19 5. X Corp. generated more than one hundred million dollars  
20 in gross revenue during the 2023 calendar year.

21 6. X is a "social media platform," as defined by § 22675(e),  
22 as it is a public internet-based application that has users in  
23 California. A substantial function of X is to connect users in  
24 order to allow users to interact socially with each other within  
25

1 the service or application. X allows users to construct public or  
2 semipublic profiles for purposes of signing into and using the  
3 application. X also allows users to create or post content viewable  
4 to others and populate a list of other users with whom an individual  
5 shares a social connection within the system.

6 **II. Burdens of AB 587**

7 7. AB 587 purports to require large social media companies,  
8 such as X Corp., to (1) post terms of service dictated by the  
9 government and include terms about how content is moderated on  
10 their platforms (the "Terms of Service Requirement") and (2)  
11 submit, on a semi-annual basis, to the California Attorney General  
12 a "terms of service report" that includes, among other things, (a)  
13 "a detailed description of content moderation practices used by  
14 the social media company for that platform"; (b) information about  
15 whether, and if so how, the social media company defines and  
16 moderates (i) hate speech or racism, (ii) extremism or  
17 radicalization, (iii) disinformation or misinformation, (iv)  
18 harassment, and (v) foreign political interference; as well as (c)  
19 information and statistics about actions taken by the social media  
20 company to moderate these categories of content (the "Terms of  
21 Service Report"). §§ 22676, 22677.

22 8. The Terms of Service Requirement in AB 587 burdens X  
23 Corp. by letting the State mandate the format and contents of X  
24 Corp.'s Terms of Service. Specifically, it requires X Corp. to,  
25

among other things:

a. Make public commitments to respond to and resolve reports of flagged content or violations of the terms of service within a specified time period. § 22676(b)(2).

b. Create and publish Terms of Service in twelve different languages - specifically, the Medi-Cal threshold languages pursuant to the Health and Safety Code § 128552, if the platform offers product features in those languages. § 22676(c). Upon information and belief this includes: (1) Arabic, (2) Armenian, (3) Cambodian, (4) Cantonese, (5) Farsi, (6) Hmong, (7) Korean, (8) Mandarin, (9) Russian, (10) Spanish, (11) Tagalog, and (12) Vietnamese. See **Exhibit 1**, which is a true and correct copy of the *Primary Language of Newly Medi-Cal Eligible Individuals*, California Department of Health Care Services, available at <https://data.chhs.ca.gov/dataset/primary-language-of-newly-medi-cal-eligible-individuals/resource/706bf0a7-9bb4-4674-9b58-917daac10d25> (last visited Oct. 6, 2023). X Corp. currently offers product features in eight of twelve applicable Medi-Cal threshold languages.

1           9. Compliance with the Terms of Service Report of AB 587  
2 would be even more burdensome - particularly in light of the  
3 quantity of content posted on X.

4           10. Our current best estimate is that, over the past 30 days,  
5 there have been on average, approximately 7,000 posts (formerly  
6 called "tweets") made on X every second. If we extrapolate from  
7 that number, that would be approximately 420,000 posts per minute,  
8 604.8 million posts per day, and almost 221 billion posts per year.

9           11. To complete the Terms of Service Report required by AB  
10 587, X Corp. would need to keep track of and categorize any and  
11 all of the content moderation decisions made as to the almost 221  
12 billion posts each year, including the millions of moderation  
13 actions taken each year by automated enforcement tools.

14           12. The Terms of Service Report required by AB 587 would  
15 force X Corp. to disclose "any existing policies intended to  
16 address" the categories of content identified by the statute (i.e.,  
17 "hate speech or racism," "extremism or radicalization,"  
18 "disinformation or misinformation," "harassment," and "foreign  
19 political interference." § 22677(a)(4)(A) and (a)(3).

20           13. This would be extremely difficult for X Corp. to do. X  
21 Corp. already has detailed descriptions of its content moderation  
22 policies available to the public online. (This is addressed in  
23 more detail below.) But X Corp.'s categories of content  
24 moderation, while comprehensive, do not align with the categories  
25

1 identified in the statute. And the listed statutory categories  
2 are controversial and difficult to define.

3 14. For example, X Corp. does not currently regulate "hate  
4 speech," "racism," or "extremism" *per se*, yet these are three  
5 categories of conduct that AB 587 forces social media companies to  
6 publicly address. § 22677(a)(3). But X Corp. regulates "violent  
7 speech," "hateful conduct," and "violent and hateful entities,"  
8 categories of content that may include content that some people  
9 might argue constitutes "hate speech," "racism," and/or  
10 "extremism." See **Exhibit 2** (Violent Speech Policy), **Exhibit 3**  
11 (Hateful Conduct Policy) and **Exhibit 5** (Violent and Hateful  
12 Entities Policy). Thus, X Corp.'s current policies do not fit  
13 within the allotted statutory categories. And it is not at all  
14 clear whether X Corp.'s "violent speech," "hateful conduct," and  
15 "violent and hateful entities" policies are "intended to address"  
16 the statutory categories, which tend to mean different things to  
17 different people.

18 15. Because the statute provides the Attorney General with  
19 discretion to impose significant civil penalties for noncompliance  
20 - of up to \$15,000 per violation per day, § 22678 - which includes  
21 a violation for material omissions or misrepresentations in the  
22 Terms of Service Report, § 22678 at (a)(2)(C), it gives the Attorney  
23 General complete discretion to determine whether he believes that  
24 it would be a violation to submit a report that did not include  
25

1 information about a policy like the "hateful conduct" policy, which  
2 does not fit neatly into the statutory categories.

3 16. The Terms of Service Report required by AB 587 would also  
4 force X Corp. to disclose "[h]ow automated content moderation  
5 systems enforce terms of service and when these systems involve  
6 human review." § 22677(a)(4)(B).

7 17. This would harm X Corp. because it would force X Corp.  
8 to disclose publicly information about its automated content  
9 moderation systems that is highly confidential and that could put  
10 X Corp. at a competitive disadvantage if made public. Moreover,  
11 disclosure of this information to the public would undermine X  
12 Corp.'s policy enforcement by providing malicious actors with  
13 information that they would likely leverage in attempts to  
14 circumvent and manipulate our policy enforcement mechanisms.  
15 Disclosure of this information is therefore likely to compromise  
16 the safety and integrity of the X platform.

17 18. Perhaps most significantly, the Terms of Service Report  
18 required by AB 587 would also force X Corp. to:

- 19 a. Collect, analyze, and generate reports detailing (i)  
20 the total number of flagged items of content; (ii) the  
21 total number of actioned items of content; (iii) the  
22 total number of actioned items of content that  
23 resulted in action taken by the social media company  
24 against the user or group of users responsible for the  
25

1 content; (iv) the total number of actioned items of  
2 content that were removed, demonetized, or  
3 deprioritized by the social media company; (v) the  
4 number of times actioned items of content were viewed  
5 by users; and (vi) the number of times actioned items  
6 of content were shared, and the number of users that  
7 viewed the content before it was actioned. §  
8 22677(a)(5)(A).

9 b. All of the information in ¶ 14(i) above must then be  
10 disaggregated into (i) the category of content;  
11 (ii) the type of content, including, but not limited  
12 to, posts, comments, messages, profiles of users, or  
13 groups of users; (iii) the type of media of the  
14 content, including, but not limited to, text, images,  
15 and videos; (iv) how the content was flagged,  
16 including, but not limited to, flagged by company  
17 employees or contractors, flagged by artificial  
18 intelligence software, flagged by community  
19 moderators, flagged by civil society partners, and  
20 flagged by users; (v) how the content was actioned,  
21 including, but not limited to, actioned by company  
22 employees or contractors, actioned by artificial  
23 intelligence software, actioned by community  
24 moderators, actioned by civil society partners, and  
25



1           actioned by users; and (vii) the number of times users  
2           appealed social media company actions taken on that  
3           platform and the number of reversals of social media  
4           company actions on appeal disaggregated by each type  
5           of action. § 22677(a)(5)(B).

6           19. The reporting requirement under AB 587 requires three  
7           reports be submitted to the Attorney General in 2024. The first  
8           report is due on January 1, 2024 and covers activity within the  
9           third quarter of 2023. § 22677(b)(2).

10          20. It would be enormously burdensome to create and  
11          categorize those records for the almost 221 billion posts made on  
12          X each year. X Corp. does not currently have the tools,  
13          infrastructure, or staff levels necessary to meet the onerous  
14          reporting requirements of AB 587. Indeed, X Corp. would need to  
15          design, build, and implement entirely new tools and workflows,  
16          including a new categorization system for moderation actions, in  
17          order to comply in good faith with the requirements of AB 587.

18          21. In the third quarter of 2023, there were approximately  
19          55 billion posts on X. X Corp. has not undertaken the burden of  
20          identifying how many of these 55 billion posts may constitute hate  
21          speech, racism, extremism, radicalization, disinformation,  
22          misinformation, harassment, or foreign political interference, and  
23          doing so would require making numerous highly controversial  
24          decisions regarding what posts fit in each of those categories.

1 Nor has X Corp. aggregated statistics on actions and appeals taken  
2 with respect to this content, as required in § 22677(a)(5).

3 22. My team's initial estimate is that implementing the  
4 infrastructure and processes necessary to comply with the  
5 requirements of AB 587 would require at least six months and involve  
6 at least thirty X Corp. employees – diverting engineering,  
7 business, and legal resources away from existing, mission-critical  
8 projects. To comply with AB 587, X Corp. would need to indefinitely  
9 commit resources to the maintenance and operation of this new  
10 compliance infrastructure.

11 23. X Corp. would need to hire new employees and/or onboard  
12 contractors in order to allocate resources to achieve good faith  
13 compliance with AB 587. Doing so would cost X Corp. hundreds of  
14 thousands or millions of dollars per year.

15 24. This does not even factor in burdens imposed by follow-  
16 up questions that will almost surely be asked by the Attorney  
17 General about compliance. Given the broad enforcement powers  
18 granted to the Attorney General under California law, AB 587 would  
19 authorize the Attorney General to issue document demands and  
20 follow-up requests to social media companies about their content  
21 moderation policies and practices to determine if they have  
22 complied with the statute in "reasonable, good faith."  
23 § 22678(a)(3).. Responding to such requests would impose  
24 significant additional burdens on X Corp. and other social media  
25

1 companies.

2 **III. The Controversial Nature of Content Moderation**

3 25. Social media content moderation is an inherently  
4 controversial undertaking. How social media companies define  
5 categories of speech that will or will not be permitted on their  
6 platforms and how they apply their rules on content moderation to  
7 particular posts are fraught with sensitive and controversial  
8 questions and are subject to rigorous debate and controversy. Put  
9 another way, deciding what content should appear on a social media  
10 platform is a question that engenders considerable debate among  
11 reasonable people about where to draw the correct proverbial line.

12 26. At X Corp., we take very seriously our responsibility to  
13 make well-thought out content moderation decisions and to make sure  
14 those decisions are accessible and understandable to our users.  
15 But in making these difficult calls, we have learned that, no  
16 matter what we do, some portion of the public will likely take  
17 issue with the way we have made these difficult judgment calls.

18 27. As the California Assembly Committee on Privacy and  
19 Consumer Protection correctly explained in analyzing AB 587,  
20 content moderation by social media companies presents those  
21 companies with a "complex dilemma" because, whatever those  
22 companies do, their decisions are subject to controversy and "both  
23 action and inaction by these companies seem[] to be equally  
24 maligned[.]" Affidavit of Joel Kurtzberg in Support of Motion for  
25

1 Preliminary Injunction, Exs. 5, 8, and 9.

2 28. This is consistent with our experience at X Corp. There  
3 is intense public debate and controversy about how to define the  
4 categories of content that should be limited on the social media  
5 platform and how to apply those categories to content on the social  
6 media platform. No matter what decisions are made, there are  
7 almost always large groups of people who disagree with them.

8 29. And the categories identified by the statute are among  
9 the most difficult to define and the most controversial to apply.  
10 While there is a general consensus that clear and extreme illegal  
11 content, such as child pornography, should not be permitted on a  
12 social media platform, there is no such consensus as to the  
13 categories identified by the statute - e.g., "hate speech, racism,  
14 extremism, misinformation, political interference, and  
15 harassment." Those categories of content are defined differently  
16 by different people and are therefore highly controversial to  
17 define and apply in practice.

18 **IV. Content Moderation at X Corp.**

19 30. X Corp. goes to great lengths to make its content  
20 moderation policies relevant and transparent. X Corp.'s rules,  
21 policies, and procedures about content moderation are all publicly  
22 available and are regularly revisited and updated based on X  
23 Corp.'s editorial judgments about what should and should not be  
24 permitted on the X platform.

1           31. Annexed hereto are true and correct copies of current X  
2 Corp. policies that concern content moderation and transparency  
3 efforts:

4           a. **Exhibit 2:** *Violent Speech Policy*, X Corp., June 2023,  
5 available at [https://help.twitter.com/en/rules-and-](https://help.twitter.com/en/rules-and-policies/violent-speech)  
6 [policies/violent-speech](https://help.twitter.com/en/rules-and-policies/violent-speech) (last visited Oct. 6, 2023);

7           b. **Exhibit 3:** *Abuse and Harassment*, X Corp., June 2023,  
8 available at [https://help.twitter.com/en/rules-and-](https://help.twitter.com/en/rules-and-policies/abusive-behavior)  
9 [policies/abusive-behavior](https://help.twitter.com/en/rules-and-policies/abusive-behavior) (last visited Oct. 6, 2023);

10          c. **Exhibit 4:** *Hateful Conduct Policy*, X Corp., Apr. 2023,  
11 available at [https://help.twitter.com/en/rules-and-](https://help.twitter.com/en/rules-and-policies/hateful-conduct-policy)  
12 [policies/hateful-conduct-policy](https://help.twitter.com/en/rules-and-policies/hateful-conduct-policy) (last visited Oct. 6,  
13 2023);

14          d. **Exhibit 5:** *Violent and Hateful Entities Policy*, X  
15 Corp., Apr. 2023, available at  
16 [https://help.twitter.com/en/rules-and-](https://help.twitter.com/en/rules-and-policies/violent-entities)  
17 [policies/violent-entities](https://help.twitter.com/en/rules-and-policies/violent-entities) (last visited Oct. 6, 2023);

18          e. **Exhibit 6:** *Abusive Profile Information*, X Corp.,  
19 available at [https://help.twitter.com/en/rules-and-](https://help.twitter.com/en/rules-and-policies/abusive-profile)  
20 [policies/abusive-profile](https://help.twitter.com/en/rules-and-policies/abusive-profile) (last visited Oct. 6, 2023);

21          f. **Exhibit 7:** *Crisis Misinformation Policy*, X Corp., Aug.  
22 2022, available at  
23 [https://help.twitter.com/en/rules-and-](https://help.twitter.com/en/rules-and-policies/crisis-misinformation)  
24 [policies/crisis-misinformation](https://help.twitter.com/en/rules-and-policies/crisis-misinformation) (last visited Oct. 6,  
25

2023);

g. **Exhibit 8:** *Synthetic and Manipulated Media Policy*, X Corp., Apr. 2023, available at <https://help.twitter.com/en/rules-and-policies/manipulated-media> (last visited Oct. 6, 2023);

h. **Exhibit 9:** *Civic Integrity Policy*, X Corp., Aug. 2023, available at <https://help.twitter.com/en/rules-and-policies/election-integrity-policy> (last visited Oct. 6, 2023).

32. Additional X Corp. rules, policies, and procedures that may address content moderation and transparency efforts are available to the public at <http://help.twitter.com/en/rules-and-policies> (last visited Oct. 6, 2023).

33. As detailed in these policies, X Corp. does moderate content that may arguably be covered by some of the controversial categories set forth in AB 587. For example, X Corp. moderates hateful conduct, crisis misinformation, violent speech, abuse and harassment, child sexual exploitation, abusive profiles, violent and hateful entities, and glorification of violence.

34. X Corp. crafted its rules, policies, and procedures with the objective of ensuring all can engage in the critical public debates surrounding these topics freely and safely.

35. In creating its rules, policies, and procedures around

1 content moderation, X Corp. takes into consideration insights  
2 gained over the past seventeen years that the platform has been  
3 operational, including feedback from trusted experts at X Corp.,  
4 our users, and the public at large.

5 36. X Corp. dedicates immense time, energy, and resources  
6 into crafting these rules, policies, and procedures and ensuring  
7 they are accessible and understandable to its users.

8 37. Through AB 587, the government is impermissibly trying  
9 to frame the content moderation debate and force social media  
10 companies like X Corp. to conform moderation practices to the  
11 government's content moderation priorities and categorizations.

12 38. By compelling speech about these controversial topics,  
13 the State is, in our view, trying to pressure X Corp. to regulate  
14 content in the way it wants by framing the public debate about  
15 content moderation. If X Corp. does not regulate the categories  
16 of content identified by the statute, then the State and/or the  
17 public may easily try to pressure X Corp. to do so by drawing  
18 attention to that fact or threatening enforcement actions unless X  
19 Corp. agrees to do so.

20 39. X Corp. would be harmed by AB 587's requirements that it  
21 make certain statements about a controversial topic against its  
22 will.

1 Dated: [REDACTED]

2 October 6, 2023

3 [REDACTED]

4 [REDACTED]

5 Sworn to before me this

6 6 day of October, 2023

7

8 See attached Certificate

9 Notary Public

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25



A notary public or other officer completing this certificate verifies only the identity of the individual who signed the document to which this certificate is attached, and not the truthfulness, accuracy, or validity of that document.

[REDACTED] } SS.

Subscribed and sworn to (or affirmed) before me on this 6 day of October, 2023, by  
 [REDACTED], proved to me on the basis of satisfactory evidence  
 to be the person(s) who appeared before me.



*[Signature]*

NOTARY'S SIGNATURE

PLACE NOTARY SEAL IN ABOVE SPACE

### OPTIONAL INFORMATION

The information below is optional. However, it may prove valuable and could prevent fraudulent attachment of this form to an unauthorized document.

#### CAPACITY CLAIMED BY SIGNER (PRINCIPAL)

- ☒ INDIVIDUAL  
☐ CORPORATE OFFICER \_\_\_\_\_ TITLE(S)  
☐ PARTNER(S)  
☐ ATTORNEY-IN-FACT  
☐ TRUSTEE(S)  
☐ GUARDIAN/CONSERVATOR  
☐ OTHER: \_\_\_\_\_

#### DESCRIPTION OF ATTACHED DOCUMENT

[REDACTED]

16  
 NUMBER OF PAGES

10/6/23  
 DATE OF DOCUMENT

ABSENT SIGNER (PRINCIPAL) IS REPRESENTING:  
 NAME OF PERSON(S) OR ENTITY(IES)

\_\_\_\_\_  
 \_\_\_\_\_

RIGHT  
 THUMBPRINT  
 OF  
 SIGNER

OTHER

